

# Do It Yourself: Generic Pseudo-labeling for Domain Adaptation in Deep Image Classification

Tsirmpas Dimitris

27th May 2024

# 1 Introduction

In the domain of deep learning, the quest for developing models that can leverage enormous sources of unlabeled data, in addition to limited labeled data, has spawned different paradigms for semi-supervised learning. In this regard, pseudo-labeling has emerged as one of the leading paradigms, promising a path toward bridging the chasm between supervised and unsupervised learning because of its simplicity, effectiveness, and adaptability across various tasks, especially computer vision tasks such as image classification.

The core of pseudo-labeling lies in the principle of iterative refinement of model predictions over unlabeled data. Given a model trained on a small labeled dataset, predictions are made on the unlabeled data. Those predictions are used to assign pseudo-labels to the unlabeled instances, effectively enlarging the labeled dataset. Further, the model is trained again on the augmented dataset in an iterative fashion, refining its performance and pseudo-labeling accuracy. It thus aims to gradually align the feature distributions of the source and target domains.

The efficacy of pseudo-labeling is conditioned on several key factors: the quality of the initially labeled data, the heuristics of pseudo-labeling, and the model architecture. These conditions have led to the inception of several algorithms, tightly coupled with specialized models able to learn domain-invariant features from both datasets. However, the level of complexity and cost to these models can be daunting, with most recent approaches necessitating multiple Deep Neural Network models (DNNs) as well as a suite of additional algorithms and training frameworks.

In this paper we propose a simple, alternative approach to training models on both unsupervised and semi-supervised datasets. Our approach seeks to minimize the training cost, only necessitating the presence of the original model, and keeping additional tuning hyper-parameters to a minimum. We evaluate the performance of our algorithm on both a small and large dataset, and extract useful insights on the weaknesses of pseudo-labeled data in the absence of the aforementioned specialized algorithms.

## 2 Related Works

## 3 Pseudo-labeling for Classification Enhancement

### 3.1 Theoretical background

Let  $D_{source} = \{x_i^s, y_i^s\}, i = 1, 2, \dots, N_{source}$  be the set containing the training samples from the source dataset, and  $D_{target} = \{x_i^t\}, i = 1, 2, \dots, N_{target}$  the equivalent from the target dataset. In domain adaptation tasks we assume that the distributions of  $D_{source}$  and  $D_{target}$  differ significantly, thus we can not simply train a classifier on  $D_{source}$  and expect satisfactory results on  $D_{target}$ .

Pseudo-labeling uses the probability estimates of a model to create labels  $y_i^t$  for the target samples  $\{x_i^t\} \in D_{target}$ . Given the estimate  $p(C|x_i^t; \Theta)$ , where  $\Theta$  the model parameters and  $C$  the class label, we assume that  $y_i^t = \text{argmax}_c p(c|x_i^t; \Theta)$ . These samples are then added to the training dataset alongside  $D_{source}$  samples. The value  $\text{max}_c p(c|x_i^t; \Theta)$  is often called the **classification confidence score** of the sample.

This strategy appears easy and intuitive but holds major caveats. If we applied the following algorithm on the entire  $D_{target}$  dataset, it is certain that the training would fail because of mislabeled samples, since, had the model's pseudo-predictions been mostly correct, there would be no incentive for the domain adaptation task in the first place. Thus, most such strategies employ an iterative approach, where only samples with a high certainty of correct classification are selected for training. This certainty is usually defined either as a high enough classification score using the same model, or via majority voting by a selection of similar models.

### 3.2 Model-agnostic Pseudo-labeling algorithm

We use a modified version of the iCAN algorithm [1]. The original iCAN algorithm was designed around a model created with the explicit purpose of learning domain-invariant features between the source and target datasets (CAN model). When the unmodified algorithm uses a generic deep model however a number of issues arise. Our approach seeks to remedy the most major of those issues.

First of all, in the original paper, the authors propose adding the loss of the fully supervised dataset  $\mathcal{L}_{source}$ , the loss of the target dataset  $\mathcal{L}_{tar}$  and the loss of their model's domain-agnostic penalty  $\mathcal{L}_{CAN}$  for each mini-batch as  $\mathcal{L} = \mathcal{L}_{source} + \mathcal{L}_{tar} + \mathcal{L}_{CAN}$ . This design decision most likely exists because of the need to have both a source and a target mini batch loaded on the network in order to calculate  $\mathcal{L}_{CAN}$ . Since this penalty does not exist in our algorithm, we instead split the training epoch into distinct backward passes for the source and target mini-batches, in order to reduce GPU VRAM requirements.

Secondly, the original iCAN algorithm selects pseudo-labeled batches for each backward pass because of the aforementioned mini-batch requirements. Since our algorithm proves much more unstable, as the underlying generic model may not learn domain-agnostic features, we instead perform pseudo-labeling once every  $N_{period}$  epochs. This mechanism ensures that our model will have acquired knowledge of the target distribution from the previously selected pseudo-labeled samples, before being asked to perform pseudo-labeling on the less-easily classifiable, unlabeled samples.

Thirdly, we do not use the "domain bias reweighing function" used by the original authors when calculating  $\mathcal{L}_{tar}$ . Aside from necessitating a second "domain" classifier, the sampling strategy we employ is inverse to the one proposed by the authors. iCAN attempts to select samples that do not fit the source distribution, in order to prevent its model from selecting target samples that are very similar to the source dataset (since they would score higher confidence scores). Our model-agnostic algorithm attempts to select samples that are closer to the source distribution and, as the model becomes more accustomed to the target distribution, slowly include samples closer to the latter. This is also the motivation behind not re-labeling pseudo-labeled samples.

The modified procedure can be found in Algorithm 1, where % is the modulo operator. The `adaptive_threshold` function is defined in Section 3.3.

---

**Algorithm 1** Modified general incremental learning algorithm

---

```

1: Train model on dataset  $D_{source}$ 
2:  $D_{pseudo} = \{\}$ 
3: for epoch do
4:   if  $epoch \% N_{period} = 0$  then
5:     Calculate accuracy on the validation source dataset
6:      $\mathcal{T} = \text{adaptive\_threshold}(\text{accuracy}, \rho)$ 
7:     for each  $d \in D_{target}$  do
8:        $label, confidence = \text{model}(d)$ 
9:       if  $confidence > \mathcal{T}$  then
10:         $D_{pseudo} = D_{pseudo} \cup \{d : label\}$ 
11:         $D_{target} = D_{target} - \{d\}$ 
12:      end if
13:    end for
14:  end if
15:   $D_{rand\_source} = \{\}$ 
16:  Select random samples from  $D_{source}$  and add to  $D_{rand\_source}$  such as  $|D_{rand\_source}| = |D_{pseudo}|$ 
17:  Train epoch on  $D_{rand\_source}$ 
18:  Train epoch on  $D_{pseudo}$ 
19: end for

```

---

### 3.3 Adaptive threshold

We use the same `adaptive_threshold` function used in the original paper, which adjusts the confidence threshold used to decide whether to pseudo-label a sample. The function is defined as  $\text{adaptive\_threshold}(acc, \rho) = \frac{1}{1 + e^{-\rho * acc}}$ , where  $acc$  is the accuracy of the classifier, and  $\rho$  a tunable hyperparameter, with  $\rho = 3$  in the original paper. We evaluate the classifier accuracy on the validation set of the source dataset.

Higher  $\rho$  values lead to a steeper decision curve, as see in Figure. The  $\rho$  parameter is thus a very convenient way of tuning the pseudo-labeling procedure; a high  $\rho$  value leads to samples chosen conservatively, ensuring that more samples are correctly labeled, while lower values lead to more samples overall being chosen, but with more incorrect labels. Since our strategy does not involve re-labeling of pseudo-labeled samples, and because of the lack of external mechanisms to detect outliers, we recommend that  $\rho \in [3, 4]$ , although the value is dependent on the datasets and underlying classifier.

Overall, the adaptive threshold function provides us with an easy-to-understand mechanism which both adapts to the current state of the model, and which can be tuned with a single hyperparameter. This hyperparameter however remains of crucial importance, and should be tuned either by the number of accepted samples, or if available, a target validation dataset which tracks the pseudo-labeling misclassifications.

### 3.4 Semi-supervised learning for learning domain-invariant features

One of the greatest challenges in the unsupervised pseudo-labeling approach is the distance between the source and target distributions, which may inhibit training. By using Algorithm 1 during experimentation we noticed three dominant patterns in the evolution of the model:

- The model selects samples liberally from  $D_{target}$ . As established in Section 3.3, this usually leads to a complete de-volution in training as the model makes increasingly more incorrect decisions, as it mislabels its own input.
- The model selects samples conservatively from a few classes. In this case, the model correctly attributes samples from a few easily distinguishable classes in  $D_{target}$  ("calculator" and "keyboard" in our experiments). It then assumes that all samples bearing features from  $D_{target}$  belong in these classes and begins mislabeling all other classes.
- The model selects samples much more conservatively but using most classes. This usually leads to a hyper-conservative model which selects samples only close to the original source distribution, and ends up disregarding a large portion of  $D_{target}$ .

In all three cases the underlying problem is that we can not, by design, expect a generic model to learn domain-invariant features.

A very simple way to partially overcome this problem is to label a small portion of  $D_{target}$  and include it during the initial pretraining phase. In other words, we pretrain the model on  $\tilde{D}_{source} = D_{source} \cup D_{labeled\_target}$ , instead of  $D_{source}$ , where  $D_{labeled\_target}$  is a stratified, labeled sample of  $D_{target}$ . This allows the model to bypass the conditions for the class imbalance problem described above, and may provide the opportunity for our model to more pick many additional samples during the first few rounds of pseudo-labeling bypassing the hyper-conservative model issue.

## References

- [1] Weichen Zhang et al. 'Collaborative and Adversarial Network for Unsupervised Domain Adaptation'. In: *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2018, pp. 3801–3809. DOI: 10.1109/CVPR.2018.00400.